

POTE Egészségügyi Szervezési Intézet

A kórházi adatfeldolgozás néhány matematikai, illetve
számítógépes problémája

Gaál Aladár

Az egészségügy egyik legsürgetőbb, megoldásra váró feladata a mindmáig korszerűtlen kórházi, rendelőintézeti adattárolás átszervezése, computerizálása. A hagyományos kartotékrendszert felszámolva, meg kell honosítanunk az elektromágneses jelrögzítést. Ez nem csupán az adatok gyorsabb hozzáférhetősége, s így a gyógyítás hatékonyságának növelése érdekében indokolt, hanem a további tudományos kutatómunka szempontjából is elengedhetetlen. Ugyanis a biológiai ismeretek komoly hányada empirikus - így a kórházak, klinikák, laboratóriumok óriási adatmennyiségéből leszűrhető tapasztalatok az alap- és alkalmazott kutatások irányvonalait is kijelölhetik.

Vitathatatlan, hogy az empirikus összefüggések alátámasztásához megkívánt nagy számú kísérlet, mérési eredmény korszerű, gyors és megbízható kiértékeléséhez nélkülözhetetlen a computerek felhasználása. Tehát a számítógépek nemcsak az adattárolás racionalizálása érdekében, hanem a bennük tárolt információk, adatok közti összefüggések felismerésének megkönnyítéséhez is szükségesek. Ezért kell nagy súlyt fektetnünk a számítógépes adatfeldolgozó módszerek kifejlesztésére.

Első lépésként az alapvető metodikai elveket, problémákat kell tisztáznunk.

A vizsgált adathalmaz változói között fennálló kapcsolatok elemzésekor két lehetőség áll előttünk:

A.) A biológiai folyamat részletes analizisével a vizsgált jelenséget több apró részfolyamatra bontjuk, s megpróbáljuk ezek számszerű kapcsolatát megállapítani. Tehát megalkotjuk a jelenség modelljét, melyben a biológiai, kémiai, fizikai, stb. hatásokat - az alkalmazott választott matematikai apparátussal mintegy összefoglaljuk. A szakismerek mélységétől függően változó tényezőket is be kell iktatnunk a modellbe. E modellt - valószínűségi változóiról - sztochasztikusnak nevezzük.

B.) Van azonban egy másik lehetőség is. Nem keressük a folyamat belső, lényegi összefüggéseit - nem a "miértekre", hanem

csak a "hogyanokra" akarunk választ kapni.

Elegendő számunkra a folyamat kezdeti (ill. bemenő) és vég (ill. kimenő) paraméterei között tapasztalható közvetlen, akár közeli-tő összefüggés ismerete is. Ugyanis sok esetben lehetetlen - a biológiai folyamatok bonyolultsága miatt - az ok-okozati összefüggéseket végigkövetni. A kauzális helyett sztochasztikus skémákkal kell dolgoznunk. Így többnyire meg kell elégednünk a jelenségre jellemző alapvető tendenciák, trendvonalak ismeretével.

1. Minőségi ismérvek:

Ha kvalitatív változókkal van dolgunk, úgy ezek előfordulásának relatív gyakoriságából tudunk következtetéseket levonni. Tekintve, hogy a vizsgált esetek általában összetettek, így nem is az egyes ismérvek, hanem ezek különböző kombinációinak gyakoriságára vagyunk kíváncsiak.

Tehát kigyűjtjük az adathalmazból, hogy hányszor fordultak elő pl. az alábbi tulajdonságok együtt:

$$X_{11} \wedge X_{21}$$

$$X_{31} \wedge X_{21}$$

$$X_{41} \wedge X_{21}$$

$$\begin{matrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{matrix}$$

$$X_{n1} \wedge X_{21}$$

$$\text{Ezen belül pl. } X_{11} \wedge X_{22} ; \quad X_{12} \wedge X_{22} ;$$

$$X_{12} \wedge X_{21} ; \quad X_{11} \wedge X_{21} \text{ stb, stb. } X_{1m} \wedge X_{2p}$$

Hasonlóképpen kigyűjthetők a három, négy, öt stb. változót tartalmazó csoportok is:

$$\text{pl. } X_{11} \wedge X_{21} \wedge X_{31} \quad \text{stb.}$$

$$X_{11} \wedge X_{22} \wedge X_{35} \wedge X_{49} \wedge \dots \wedge X_{nm} \quad \text{stb, stb.}$$

(Megjegyzés: az első index a tulajdonságra, a második annak állapotára utal.)

Kombinatorikai ismeretek birtokában belátható, hogy - a változók különböző állapotából adódó eshetőségeket nem is tekintve, - az egyes változó-csoportok lehetséges előfordulási száma n db változó esetén:

$$\sum_{k=0}^n \binom{n}{k} = 2^n$$

Ez - láthatóan - a figyelemmel kísért változók számával exponenciálisan nő.

Ezen belül pl. egy adott $k=3$ -as változócsoporthoz további p, q, r eshetőség van, ha a változók állapotának száma p, q és r .

Miután a gyakorlatban - egy-egy nagyobb felmérésnél - n értéke százas nagyságrendben mozog, a tulajdonságcsoporthoz előfordulási gyakoriságait csakis számítógépekkel érdemes kigyűjteni.

A szakembernek csak a válogatás főbb irányait kell megjelölni, az idő- és energia-pazarló rutinmunkát a gép végzi el. Természetesen a legkorszerűbb gépi adatfeldolgozás sem nélkülözheti a gondolkodó embert! Erre nemcsak az adatok betáplálását megelőző szelektálásnál, hanem az eredmények értékelésénél is szükség van. Fantos, hogy a nem vizsgált változók hatását elimináljuk!

Miután a biológiai vizsgálatoknál általában nincs módunk a jelenségeket tetszés szerint befolyásolni, irányítani, így nem tudjuk lejátszani, modellezni úgy a folyamatokat, hogy az egyes tényezőket rendre kiküszöbölve, az egyedi hatásokat megállapíthassuk. Ezért úgy kell válogatnunk az adatokból, hogy a nem vizsgált változók értékei állandók legyenek (ill. egy elfogadhatóan szűk intervallumon belül mozogjanak), hogy hatásuk ne érvényesüljön: ne fedhessék el a vizsgált kapcsolatot.

Pl. ha az X_1 (bemeneti) változónak az Y (kimeneti) ordinátóra gyakorolt parciális hatását vizsgáljuk, úgy a többi X_2, X_3, \dots, X_n változónak konstansnak kell lenniük. Ellenkező esetben ugyanis nem lehetünk biztosak abban, hogy Y megváltozásait csak X_1 megváltozásai okozzák - hiszen a "nem figyelt" $X_2, X_3, X_4, \dots, X_n$ tényezők is hatnak, akaratunktól függetlenül, objektívek.

Ha nem így járunk el, úgy akár a valósággal homlokegyenest ellenkező eredmények is adódhatnak!

Igy a rendelkezésünkre álló adathalmazt valamennyi változó szerint, számtalan szempontból analizálhatjuk, - mely gépi úton könnyen és gyorsan elvégezhető.

Természetesen az értékelés nagy körültekintést és megfelelő szakmai ismereteket kíván meg. Az egyes tulajdonságok, ismérvek együttes előfordulásának gyakoriságából nem lehet egyértelműen ok-okozati kapcsolatra következtetni. Lehet ritkán előforduló faktorok között is összefüggés - s megfordítva is áll: a gyakran együttjáró tulajdonságok sem jelentenek mindig kauzalitást.

II. Mennyiségi ismérvek:

Amennyiben adataink kvantitatívek, úgy a regresszió- ill. faktoranalízis segítségével kereshetjük a vizsgált tényezők között fennálló kapcsolatokat.

A.) A regressziószámítás lényegében a mért változók valamennyi adatához legjobban illeszkedő függvény együtthatóit határozza meg. (Gauss-féle legkisebb négyzetek elve.) Lineáris összefüggés esetén az

$$Y = a_0 + a_1 X_1 + a_2 X_2 + \dots + a_n X_n$$

egyenlet együtthatóira teljesülnie kell az alábbiak:

$$\sum_{i=1}^m \left[Y_i - (a_0 + a_1 X_{1i} + a_2 X_{2i} + \dots + a_n X_{ni}) \right]^2 = \varnothing \stackrel{!}{=} \text{Minimum}$$

/m = esetszám/

Ebből:

$$\frac{\partial \varnothing}{\partial a_0} \stackrel{!}{=} 0$$

$$\frac{\partial \varnothing}{\partial a_1} \stackrel{!}{=} 0$$

$$\vdots$$

$$\frac{\partial \varnothing}{\partial a_n} \stackrel{!}{=} 0$$

A többváltozós szélsőértékszámítás fenti egyenleteiből a keresett $a_0, a_1, a_2, \dots, a_n$ regressziós együtthatók adódnak.

Ezek, mint az egyes változók súlyozó együtthatói, megmutatják, hogy a kérdéses tényezők milyen mértékben és irányban befolyásolják a vizsgált folyamat eredményét, külön-külön.

A regressziós egyenletből tehát választ kapunk a jelenséget befolyásoló be- és kimenő változók, azaz az okok és okozatok közti közvetlen - a részeffektusokat nem tükröző - kapcsolat jellegére. Nem ismervén, illetve nem vizsgálván a folyamatot befolyásoló valamennyi tényező hatását, e kapcsolat természetesen csak mint tendencia juthat érvényre. Így a regressziós összefüggés segítségével - egy jövőbeni mérés, illetve kísérlet eredményére csak egy megadott valószínűséghez rendelhető bizonytalansági közön belül következtethetünk. (Konfidencia-intervallum).

Ha a változók mérési tartományán kívül becslést alkalmazunk, úgy következtetéseinknek szakmailag is megalapozottnak kell lenniük. Ugyanis maga a statisztikai elemzés az extrapolációra nem jogosít fel.

Míg a regressziószámítás a változók közti kapcsolat jellegét, a korrelációszámítás annak szorosságát adja meg.

A többszörös korrelációs együttható az összes független változót egyként kezelve - méri azok és a függő változó közti kapcsolat mértékét.

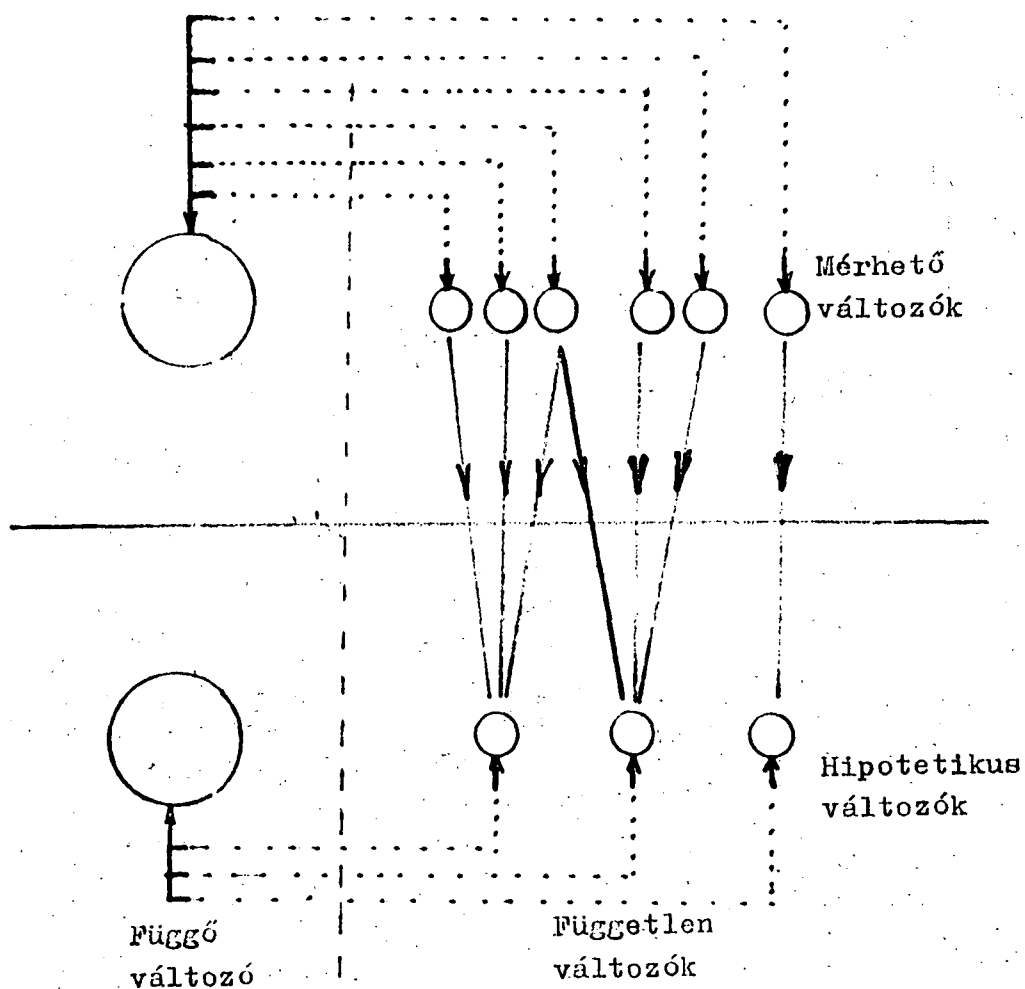
A parciális korrelációs együttható viszont az egyes független változók és a függő változók közti korrelációt jelzi. A többi független változót állandó értéken tartva, azok hatása nem érvényesül.

Ezzel szemben a totális korrelációs együtthatók a két kiszemelt változó közti közvetlen és közvetett hatásokat együttesen mérik. Ugyanis a kiválasztott két változó a többi változóval kapcsolatban van, - így e mérőszámban az azokon keresztül érvényesülő hatásuk is kifejeződik.

B.) Faktoranalízis:

Szemben a regresszióanalízissel, mely a mérhető változók kapcsolatát elemezte, - a faktoranalízis már "mélyebb rétegekben" keresi az okozati összefüggéseket.

Sematikusan:



Ugyanis a vizsgált (mért) "független" változók valójában nem minden esetben függetlenek egymástól, mert ezek - illetve egy részük - számunkra még ismeretlen közös faktorok: ún. "hipotetikus változók" függvényei. A mért változók ezeken a feltételezett közös faktorokon keresztül kerülnek egymással függőségbe.

Igy az egyes $X_1, X_2, X_3, \dots, X_n$ mért változók valamilyen $K_1, K_2, K_3, \dots, K_m$ közös faktorok lineáris kombinációjaként előállíthatók:

$$X_1 = a_{11} K_1 + a_{12} K_2 + \dots + a_{1m} K_m$$

$$X_2 = a_{21} K_1 + a_{22} K_2 + \dots + a_{2m} K_m$$

$$\vdots$$

$$X_n = a_{n1} K_1 + a_{n2} K_2 + \dots + a_{nm} K_m$$

De ha a mért változók előállíthatóak a közös faktorok felhasználásával, akkor megfordítva is igaz: az egyes közös faktorok is kell, hogy előállíthatóak legyenek a ténylegesen megfigyelt változók segítségével. Igy pl. a K_1 -re írható:

$$K_1 = \alpha_1 X_1 + \alpha_2 X_2 + \dots + \alpha_n X_n$$

Tehát a közös faktorok - melyek jelenlétére a mért változók függőségéből következettünk, - számszerűsíthetők, a mérési értékekből - az α_i együtthatók ismeretében megadhatóak.

Vizsont számtalan esetben nehézséget okoz ezek biológiai értelmezése! Ugyanis ha egy-egy hipotetikus változó a legkülönbözőbb jellegű mért változók lineáris kombinációjaként adódik, - úgy nehéz annak biológiai tartalmat tulajdonítani.

Amennyiben így egy újabb kapcsolatot ismerünk fel, úgy hasznos lehet e módszer, egyébként pedig - nem mondván lényegesen többet a változók kapcsolatáról, mint a regresszióanalízis - megfontolandó, hogy indokolt-e ahelyett a faktoranalízis sokkalta bonyolultabb apparátusát belevonni vizsgálatainkba. Ugyanis ha a megfigyelt változók standardizált (z) értékeire (amikor is azok átlaga zérus, szórásuk pedig egy) általánosán:

$$\underline{z} = \underline{A} \cdot \underline{f} = \underline{A}_k \cdot \underline{k} + \underline{A}_u \cdot \underline{u}$$

ahol:

\underline{z} = a standardizált változók oszlopvektora

\underline{k} = a közös faktorok oszlopvektora

\underline{u} = az un. egyedi faktorok oszlopvektora

\underline{A} = együttható-mátrix

akkor ebből - részletesen kiírva - speciális alakú regressziós egyenletet kapunk.

Egy adott z_i változóra ez:

$$z_i = a_{i1}K_1 + a_{i2}K_2 + \dots + a_{im}K_m + a_{ij}U_j$$

Az ebben szereplő független (K) változók nem mérhetőek, így az a_{ij} faktorsúlyok sem határozhatóak meg a regressziós elemzés szokásos módszereivel - csak más, lényegesen bonyolultabb uton.

A részletek mellőzésével, az un. "alapvető faktorok módszere" (Principal Factor Solution) az alábbi sajátérték - sajátvektor problémára vezet:

$$(\underline{R}_h - \lambda_1 \underline{E}) \cdot \underline{a}_1 = \underline{0}$$

azaz:

$$\underline{R}_h \cdot \underline{a}_1 = \lambda_1 \cdot \underline{a}_1 \quad \underline{R}_h = \frac{1}{N} \underline{Z}' \underline{Z}^*$$

ahol:

\underline{E} = egységmátrix

N = az elemszám

\underline{R}_h = a megfigyelések redukált korrelációs mátrixa

(Megjegyzés: itt az un. egyedi faktoroktól eltekintettünk: a \underline{Z}' erre utal.)

\underline{a}_1 tehát az \underline{R}_h redukált korrelációs mátrix egy sajátvektora.

Egyéb megszorításokat is figyelembevéve az \underline{a}_1 vektort az \underline{R}_h legnagyobb saját-értékéhez tartozó saját-vektorok halmazából kell kiválasztanunk.

Ez éppen egy $\sqrt{\lambda_1}$ hosszúságú sajátvektor.

Az \underline{A} mátrix második \underline{a}_2 oszlopának meghatározása az \underline{a}_1 meghatározásához hasonlóképpen történik, azzal a különbséggel, hogy \underline{R}_h szerepét az

$$\underline{R}_1 = \underline{R}_h - \underline{a}_1 \underline{a}_1^*$$

un. első reziduális korrelációs mátrix veszi át. A továbbiakban is hasonlóképpen járunk el.

Az így nyert $\underline{a}_1, \underline{a}_2, \dots, \underline{a}_m$ vektorok páronként ortogonálisak, és az eljárás során adódó λ sajátértékekre fennáll:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$$

Az utolsó \underline{R}_m reziduális mátrix értéke zérus, így a fenti eljárás teljes mértékben reprodukálja a redukált korrelációs mátrixot.

A vázolt módszerrel tehát az \underline{a}_{ij} faktorsúlyok - s így - átrendezés után az α_i keresett együtthatók is adódnak. Azaz a mért változók segítségével ez utóbbiak ismeretében a K_i közös faktorok meghatározhatók.

Hogy az ismertett módszerek közül mikor melyiket célszerűbb használni, azt esetenként kell eldönteni. Intézetünkben (Pécsi Orvostudományi Egyetem Egészségügyi Szervezéstudományi Intézet Számítástechnikai Csoport) az elmondottak figyelembevételével dolgoztuk ki a számítógépes adatfeldolgozó programokat.

E programokkal vizsgáljuk szociális körülmények, társadalmi adottságok és a különféle betegségek kapcsolatát.

Segítségükkel elemezzük a koraszülést befolyásoló tényezőket is. Ugyanakkor lehetővé tesszük a klinikákon, kórházakban felgyűlemlett hatalmas adathalmaz korszerű, gyors feldolgozását, áttekintését: kezelési eljárások, gyógyszerhatások és mellékhatások, laboratóriumi eredmények, diagnózisok stb. közti összefüggések felismerését.

I r o d a l o m

- 1.) Ezekiel, M. - Fox, K.: Korreláció és regresszió analízis
(lineáris és nem lineáris módszerek)
Közgazdasági és Jogi Kiadó Budapest, 1970.
- 2.) Prékopa András: Valószínűségelmélet. (Műszaki alkalmazásokkal.)
Műszaki Könyvkiadó Budapest, 1972.
- 3.) Vincze István: Matematikai statisztika. (Ipari alkalmazásokkal.)
Műszaki Könyvkiadó Budapest, 1968.
- 4.) Vita László: A faktoranalízis közgazdasági alkalmazásának lehetőségeiről.
Sigma 1970. III. évf. 2.sz. 127-152 old.